

Lab report: Improving a naive matrix product implementation

Nicolas Daix-Moreux

[Link to github repo](#)

Part 1: Naive Matrix Product Performance & Scaling Study

System Specifications

Component	Description
CPU	Intel® Core™ i5-1035G1 @ 1.00GHz (4 cores, 8 threads)
Max Clock Speed	3.60 GHz
Architecture	x86_64 (64-bit support)
Cache	L1: 192 KiB (data) + 128 KiB (instruction) L2: 2 MiB L3: 6 MiB
Virtualization	VT-x
NUMA Nodes	1 node with CPUs 0–7
Memory	7.8 GiB total online memory 7.3 GiB usable 2.0 GiB swap
GPU	Intel Iris Plus Graphics G1 (Ice Lake)
Memory Blocks	128 MiB block size; removable memory blocks supported
OS	(Assumed from shell) Ubuntu or Debian-based Linux (64-bit)

Introduction

This document reports on the performance of a naive matrix multiplication using [Kokkos](#) in C++. We explore both **runtime** and **GFLOP/s**, conduct a **strong scaling study**, and compare two data layouts: **LayoutLeft** vs **LayoutRight**.

```
template <class AMatrixType, class BMatrixType, class CMatrixType>
auto matrix_product(double alpha, AMatrixType const& A, BMatrixType const&
B, double beta, CMatrixType& C) -> void {
```

```
static_assert(
    AMatrixType::rank() == 2 && BMatrixType::rank() == 2 &&
    CMatrixType::rank() == 2, "Views must be of rank 2"
);
assert(A.extent(0) == C.extent(0));
assert(B.extent(1) == C.extent(1));
assert(A.extent(1) == B.extent(0));

Kokkos::parallel_for(
    "dgemm_kernel",
    A.extent(0),
    KOKKOS_LAMBDA(int i) {
        for (int j = 0; j < int(B.extent(1)); ++j) {
            double acc = 0.0;
            for (int k = 0; k < int(A.extent(1)); ++k) {
                acc += alpha * A(i, k) * B(k, j);
            }
            C(i, j) *= beta + acc;
        }
    }
);
}
```

1. Measuring Performance for the naive matrix product

We run `strong_scaling.py` for the following tests. We measure three primary metrics:

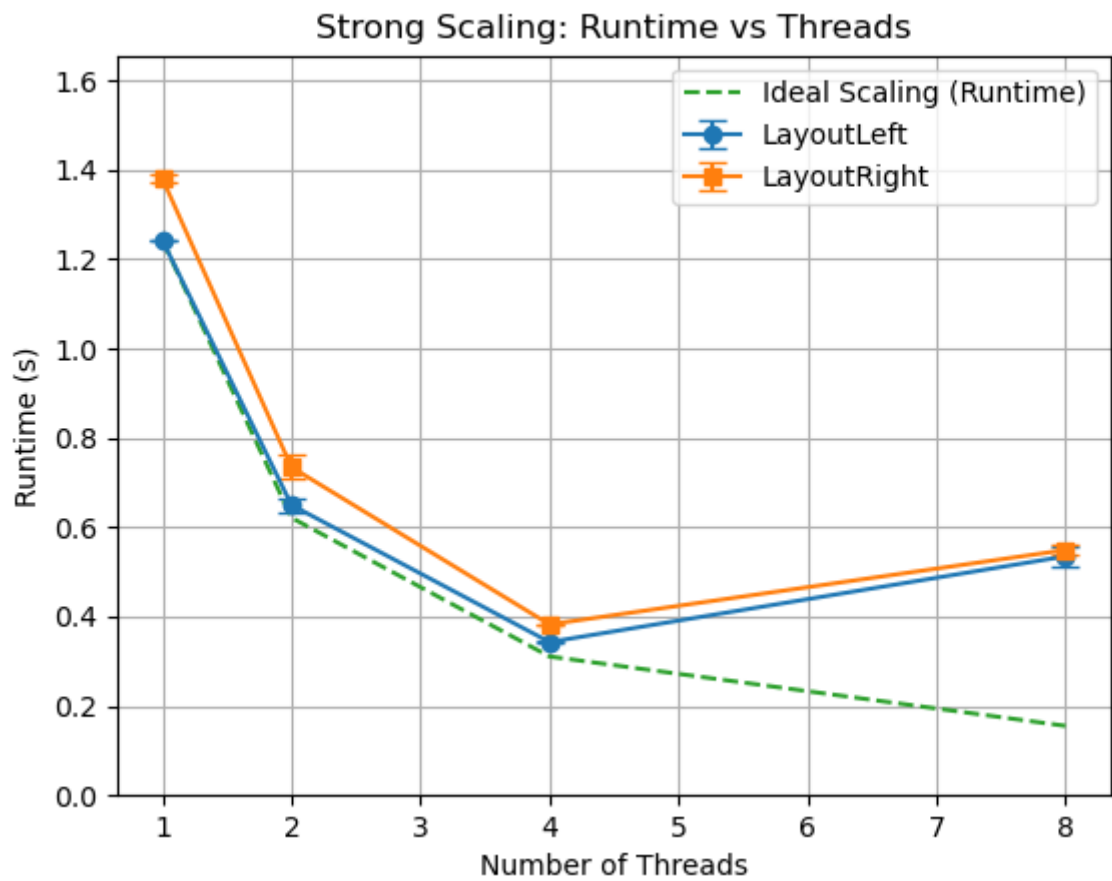
- **Runtime (s):** Wall-clock time for the entire multiplication.
- **GFLOP/s** = $2 * M * N * K * 10^9$
- **Speedup:** Runtime gains achieved relative to thread count

Strong Scaling Results

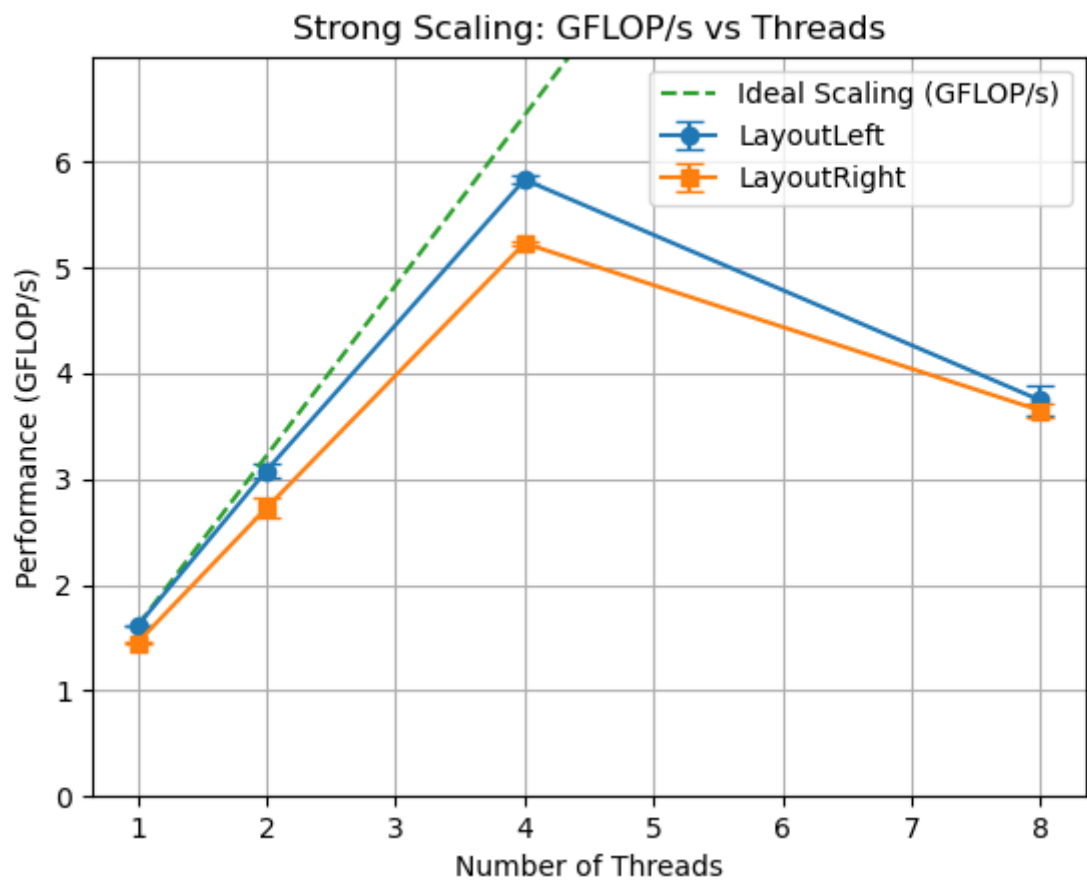
Threads	Layout	Runtime (s)	GFLOP/s
1	Left	1.2424 ± 0.0019	1.61 ± 0.00
2	Left	0.6497 ± 0.0156	3.08 ± 0.07
4	Left	0.3429 ± 0.0020	5.83 ± 0.03
8	Left	0.5349 ± 0.0217	3.74 ± 0.15
1	Right	1.3801 ± 0.0092	1.45 ± 0.01
2	Right	0.7344 ± 0.0268	2.73 ± 0.10
4	Right	0.3824 ± 0.0010	5.23 ± 0.01
8	Right	0.5490 ± 0.0114	3.65 ± 0.07

Plots are located in the `plots_strong_scaling` folder.

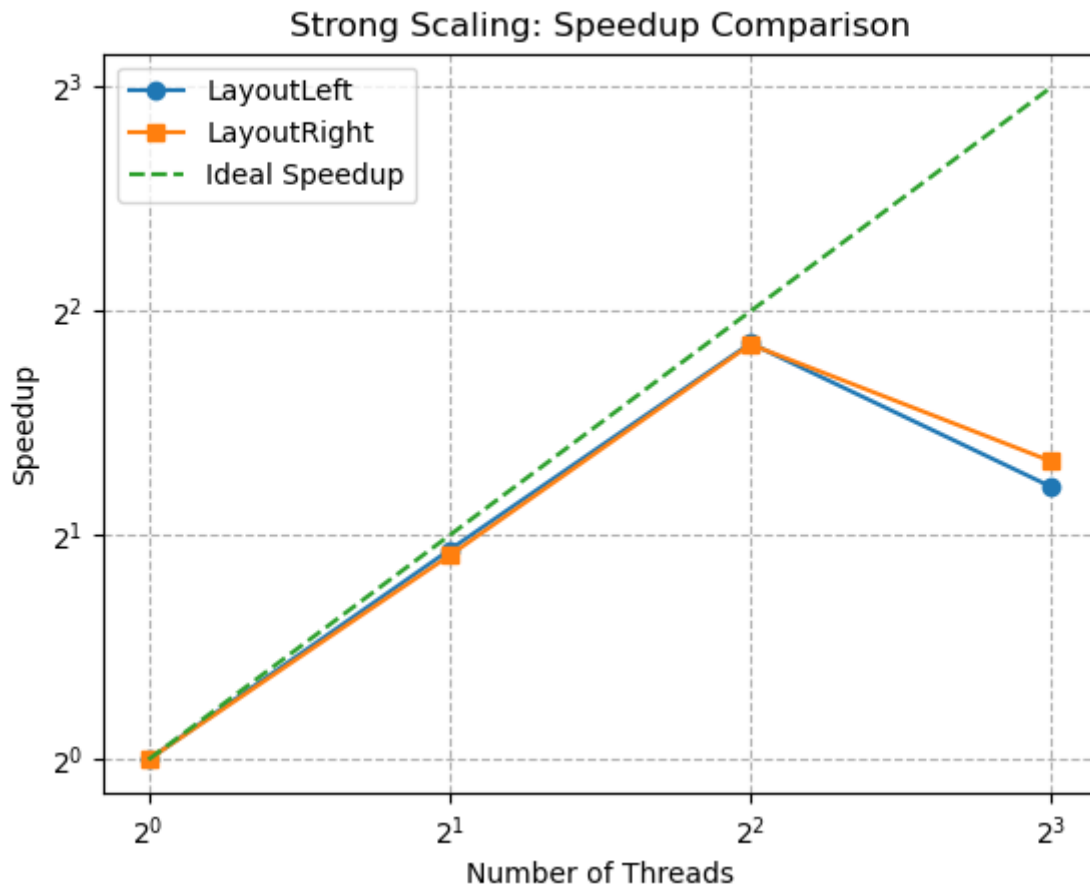
1. Runtime Comparison



2. GFLOP/s Comparison



3. Speedup Comparison



Analysis of Kokkos Layout Performance

Layout Comparison

The observed results clearly demonstrate a difference in performance between the two memory layouts used by Kokkos:

- **Kokkos::LayoutLeft** consistently outperforms **Kokkos::LayoutRight** in every tested configuration. The advantage is especially clear at lower thread counts, e.g., at a single thread (1.61 vs. 1.45 GFLOP/s) and at two threads (3.08 vs. 2.73 GFLOP/s).
- At higher thread counts (4 and 8 threads), LayoutLeft continues to outperform LayoutRight, achieving a higher peak performance (5.83 GFLOP/s for LayoutLeft vs. 5.23 GFLOP/s for LayoutRight).

Reasons for Performance Difference

The difference in performance arises from the underlying memory access patterns:

- **Kokkos::LayoutLeft** is column-major, meaning consecutive memory elements in a column are stored contiguously. This layout typically benefits CPU cache performance when accessing elements column-wise, as is commonly done in optimized matrix multiplications.
- **Kokkos::LayoutRight** is row-major (the default C++ array layout), which arranges consecutive elements of a row contiguously. While intuitive in typical C++ contexts, it can lead to less efficient memory access patterns for certain linear algebra operations.

Due to better alignment with memory access patterns in the matrix multiplication kernel, LayoutLeft achieves better cache utilization, significantly reducing cache misses and thus improving performance.

Strong Scaling Analysis

Runtime

- For **LayoutLeft**, runtime initially scales very well from 1 to 4 threads (1.24s → 0.34s). However, adding more threads (from 4 to 8) worsens the runtime (0.34s → 0.53s), indicating that optimal parallel efficiency is reached at 4 threads.
- For **LayoutRight**, the trend is similar: runtime improves significantly from 1 to 4 threads (1.38s → 0.38s), and then degrades at 8 threads (0.38s → 0.55s).

GFLOP/s

- **Ideal scaling** (doubling threads doubles performance) occurs from 1 to 2 threads in both layouts, and nearly ideal from 2 to 4 threads.
- At 8 threads, both configurations exhibit a significant decrease in performance, strongly suggesting thread contention or that hyper-threading is not beneficial for the tested kernel.

Speedup

- **LayoutLeft** achieves near-ideal linear speedup from 1 to 4 threads, with a maximum speedup of ~3.6x.
- **LayoutRight** also shows near-linear scaling up to 4 threads but with slightly lower peak efficiency (~3.6x speedup).
- At 8 threads, efficiency drops significantly for both layouts.

Conclusion

The optimal configuration identified through this analysis is clearly **Kokkos : LayoutLeft at 4 threads**, yielding the highest performance and most effective memory utilization. Running the matrix product with 8 threads results in a performance drop, this indicates potential resource contention or hyper-threading inefficiencies beyond the number of physical CPU cores available (which are typically 4 physical cores for an Intel Core i5-1035G1). Future optimizations might focus on mitigating thread contention and optimizing cache usage further at higher thread counts.

Part 2: Cache Blocking

Introduction

We are now going to optimize the previous naive matrix product by utilizing cache blocking and adjusting block size to try and squeeze out as much performance as possible out of the

All runs are now done on the **LayoutLeft** version of the code, **OMP_PLACES="cores"**, **OMP_PROC_BIND=true**, **OMP_NUM_THREADS=4** if not specified otherwise.

Q1 Perf profiling

Quick Introduction to the Perf Tool

perf is a powerful Linux profiling tool used to analyze performance aspects of software, particularly CPU performance. It accesses CPU hardware counters that measure various events, such as:

- **Cache usage** (hits and misses at different cache levels).
- **Branch predictions** (correct/incorrect predictions).
- **Instructions executed**, cycles consumed, and more.

Key points about **perf**:

- **Hardware Counters:** **perf** uses built-in CPU performance counters to track events.
- **Events:** You can monitor specific events (**cache-misses**, **cache-references**, **L1-dcache-loads**, etc.).
- **Overhead:** **perf** incurs minimal performance overhead, making it suitable for production-level profiling.

Analysis of **perf** Results

We run the following command:

```
perf stat -e L1-dcache-loads,L1-dcache-load-misses,LLC-loads,LLC-load-misses \
./top.matrix_product 1000 1000 1000
```

Metric	Value	Interpretation
Runtime	0.365 s	Reasonable overall performance
Performance	5.48 GFLOP/s	Moderate computational throughput
L1 Cache Loads	1,599,912,257	
L1 Cache Missses	1,156,202,708	
L1 Cache Miss Rate	72.27% (very high)	Poor L1 cache locality; memory access issues
LLC Cache loads	5,374,711	
LLC Cache misses	3,858,748	
LLC Miss Rate	71.79% (very high)	Poor LLC reuse; large memory-bound penalty

2. Cache Performance Analysis

L1 Cache Analysis:

This is a very high L1 miss rate (ideally should be much lower, typically below 5-10%). This means about **3 out of every 4 memory accesses miss the fastest cache (L1)**, resulting in higher latency accesses and significantly impacting performance.

This is typically due to:

- Poor memory access patterns (e.g., not exploiting spatial locality).
- Large working sets not fitting within the small L1 cache (only 48 KiB per core in your CPU).

Last-Level Cache (LLC/L3) Analysis:

This is also a **high LLC miss rate**, meaning when data misses in L1 and then eventually misses in the LLC, it must fetch data directly from main memory (RAM), significantly increasing latency.

High LLC miss rates typically indicate:

- Data access patterns that exceed the cache size (your CPU's LLC is 6 MiB).
- Insufficient data reuse (low temporal locality).

Interpretation

The provided **perf** results highlight clear performance bottlenecks in memory access:

- **High L1 Miss Rate:**
 - Suggests the blocked or loop-ordered matrix multiplication is not optimized for spatial locality, causing frequent misses even at the fastest cache level.
 - This typically occurs if loops access memory inefficiently (e.g., accessing matrices in a sub-optimal pattern, like accessing columns instead of rows, depending on memory layout).
- **High LLC Miss Rate:**
 - Indicates the working set of your problem does not fit within the CPU's largest cache (L3 cache).
 - Frequent misses at this level dramatically slow performance because memory fetches must occur from main memory, typically ~100 times slower than L1 cache accesses.

Q2 Implementing a cache blocked version

To significantly enhance the performance of the naive matrix multiplication, I implemented a cache-blocked version of the algorithm. The goal of this approach is to improve both spatial and temporal data locality, reducing cache misses and thus increasing computational efficiency.

Cache-blocked Implementation

The provided cache-blocked function splits matrices into smaller sub-blocks defined by block sizes **BM**, **BN**, and **BK**. These sub-blocks are specifically chosen to fit within the CPU cache hierarchy (primarily L1 and L3 caches). Each sub-block of matrices A, B, and C is loaded into local arrays (**local_C**) within the loop, which greatly enhances data reuse:

```
template <class AMatrixType, class BMatrixType, class CMatrixType>
void matrix_product_block(double alpha, AMatrixType const& A, BMatrixType
const& B, double beta, CMatrixType& C, int BM, int BN, int BK) {
    static_assert(
        AMatrixType::rank() == 2 && BMatrixType::rank() == 2 &&
```

```

CMatrixType::rank() == 2,
    "Views must be of rank 2"
);
assert(A.extent(0) == C.extent(0));
assert(B.extent(1) == C.extent(1));
assert(A.extent(1) == B.extent(0));

// Define block sizes (tune these based on the CPU cache)
///Default values for BM,BN,BK is 32, defined in the main function

int BLOCK_I = BM;
int BLOCK_J = BN;
int BLOCK_K = BK;

using execution_space = typename CMatrixType::execution_space;
using policy_t = Kokkos::MDRangePolicy<
    execution_space,
    Kokkos::Rank<2, Kokkos::Iterate::Default, Kokkos::Iterate::Default>,
    Kokkos::IndexType<int>
>;

// Number of blocks in each dimension
int num_blocks_i = (C.extent(0) + BLOCK_I - 1) / BLOCK_I;
int num_blocks_j = (C.extent(1) + BLOCK_J - 1) / BLOCK_J;

Kokkos::parallel_for("dgemm_blocked", policy_t({0, 0}, {num_blocks_i,
num_blocks_j}),
KOKKOS_LAMBDA(const int bi, const int bj) {
    // Define the block boundaries
    const int start_i = bi * BLOCK_I;
    const int end_i = (bi + 1) * BLOCK_I < C.extent(0) ? (bi + 1) * BLOCK_I
: C.extent(0);
    const int start_j = bj * BLOCK_J;
    const int end_j = (bj + 1) * BLOCK_J < C.extent(1) ? (bj + 1) * BLOCK_J
: C.extent(1);

    // Temporary local storage for the C block
    double local_C[BLOCK_I][BLOCK_J] = {};

    // Initialize local_C with beta * C
    for (int i = start_i; i < end_i; ++i) {
        for (int j = start_j; j < end_j; ++j) {
            local_C[i - start_i][j - start_j] = beta * C(i, j);
        }
    }

    // Loop over K in blocks
    const int num_blocks_k = (A.extent(1) + BLOCK_K - 1) / BLOCK_K;
    for (int bk = 0; bk < num_blocks_k; ++bk) {
        const int start_k = bk * BLOCK_K;
        const int end_k = (bk + 1) * BLOCK_K < A.extent(1) ? (bk + 1) *
BLOCK_K : A.extent(1);

        // Compute contribution from this K block

```



```
    for (int k = start_k; k < end_k; ++k) {
        for (int i = start_i; i < end_i; ++i) {
            const double a = alpha * A(i, k);
            for (int j = start_j; j < end_j; ++j) {
                local_C[i - start_i][j - start_j] += a * B(k, j);
            }
        }
    }

    // Write back to global C
    for (int i = start_i; i < end_i; ++i) {
        for (int j = start_j; j < end_j; ++j) {
            C(i, j) = local_C[i - start_i][j - start_j];
        }
    }
});
}
```

Performance Analysis Using Perf

The cache-blocked implementation was profiled using **perf** to measure cache performance for matrices of size **1000×1000×1000**. The key performance metrics obtained were:

Metric	Naive Version	Cache-blocked Version	Improvement
Runtime (seconds)	0.3649 s	0.1397 s	2.61×
Performance (GFLOP/s)	5.48 GFLOP/s	14.31 GFLOP/s	2.61×
L1 Cache Miss Rate	72.27%	0.55%	131×
LLC (L3) Miss Rate	71.79%	9.22%	7.79×

The introduction of the cache-blocked implementation resulted in remarkable performance enhancements:

- **Speedup:** The cache-blocked version achieves a **2.61×** improvement in runtime and performance.
- **Cache Efficiency:** It notably improves L1 and L3 cache efficiency, with L1 miss rates dropping by **131×** and L3 misses by **7.79×**.

Explanation of Performance Improvements

The observed performance gains and reductions in cache misses are due to the following factors:

1. Local Buffering for Temporal and Spatial Locality:

- **Local Copy:**
For each block of **C**, the kernel allocates a small, stack-based two-dimensional array **local_C** of size **BM*BN**). This local array is used to hold the partial results of the block multiplication. Initially, it is loaded with **beta*C**, ensuring that the final update preserves the contribution from the original **C** (scaled by **beta**).

- **Blocking Over K :**

Instead of processing the entire K dimension at once, the kernel divides K into smaller chunks of size BK . For each such sub-block, the corresponding sub-tiles of A and B are used to update $local_C$. Because these updates occur in the fast, on-chip (or register) memory, the data can be reused multiple times with very low latency.

- **Improved Locality:**

By keeping each $BM \times BN$ block in local memory, the algorithm reduces the number of global memory accesses for C during the innermost loops. This leads to fewer L1 cache misses and better spatial locality, as the local array is stored contiguously in memory.

2. Parallelization via MDRangePolicy:

- **2D Tiling:**

The kernel uses a 2D MDRangePolicy to divide the work over the i and j dimensions (the block indices of C). Each iteration of the policy handles one block (tile) of C . Because these blocks are processed independently, they can be executed in parallel without inter-dependency.

- **Reduced Synchronization Overhead:**

Once the local block computation is complete, the result is written back to the global memory. This design minimizes the need for inter-thread synchronization, as each block's computation is self-contained.

- **Load Balancing:**

Partitioning work in this way provides a natural granularity for parallel execution. However, the number of blocks per matrix (especially for smaller matrices) may limit the scaling. For larger matrices, however, many more blocks are available, which helps in distributing the work evenly among available threads.

3. Impact on Cache Miss Rates and Performance:

- **Spatial Locality:**

The contiguous allocation of $local_C$ ensures that once a block of data is brought into cache, multiple operations are performed on it, thereby reducing the effective cache miss rate.

- **Temporal Locality:**

Each block is reused across all iterations over the K sub-blocks before the final write-back to C , which helps to amortize the cost of loading data from main memory.

- **Overall Throughput:**

The experimental results show that, by carefully tuning the block dimensions, significant performance improvements can be achieved (e.g., for the $1000 \times 1000 \times 1000$ matrix, a configuration such as $32 \times 32 \times 32$ yields a runtime of 0.14 s and 14.3 GFLOP/s compared to far lower throughput for less-optimal configurations).

Q3 Verification of cache blocked function

Added a `top.matrix_verify` in `/src` that runs the code sequentially after running with the blocked version, and allows for up to 0.01 relative error. This error stems from the rearrangement of the operations performed at compile time and runtime. On average, the error is of about 0.06 %.

Q4 Tuning the Block Dimensions for Cache-blocked Matrix Multiplication

The objective of this experiment was to identify the optimal block dimensions (**BM**, **BN**, **BK**) for my cache-blocked matrix multiplication implementation, maximizing computational performance and cache efficiency. Using the Perf profiling tool, I systematically tested multiple configurations to measure runtime, computational performance (GFLOP/s), and cache miss rates (L1 and LLC).

Experimental Results

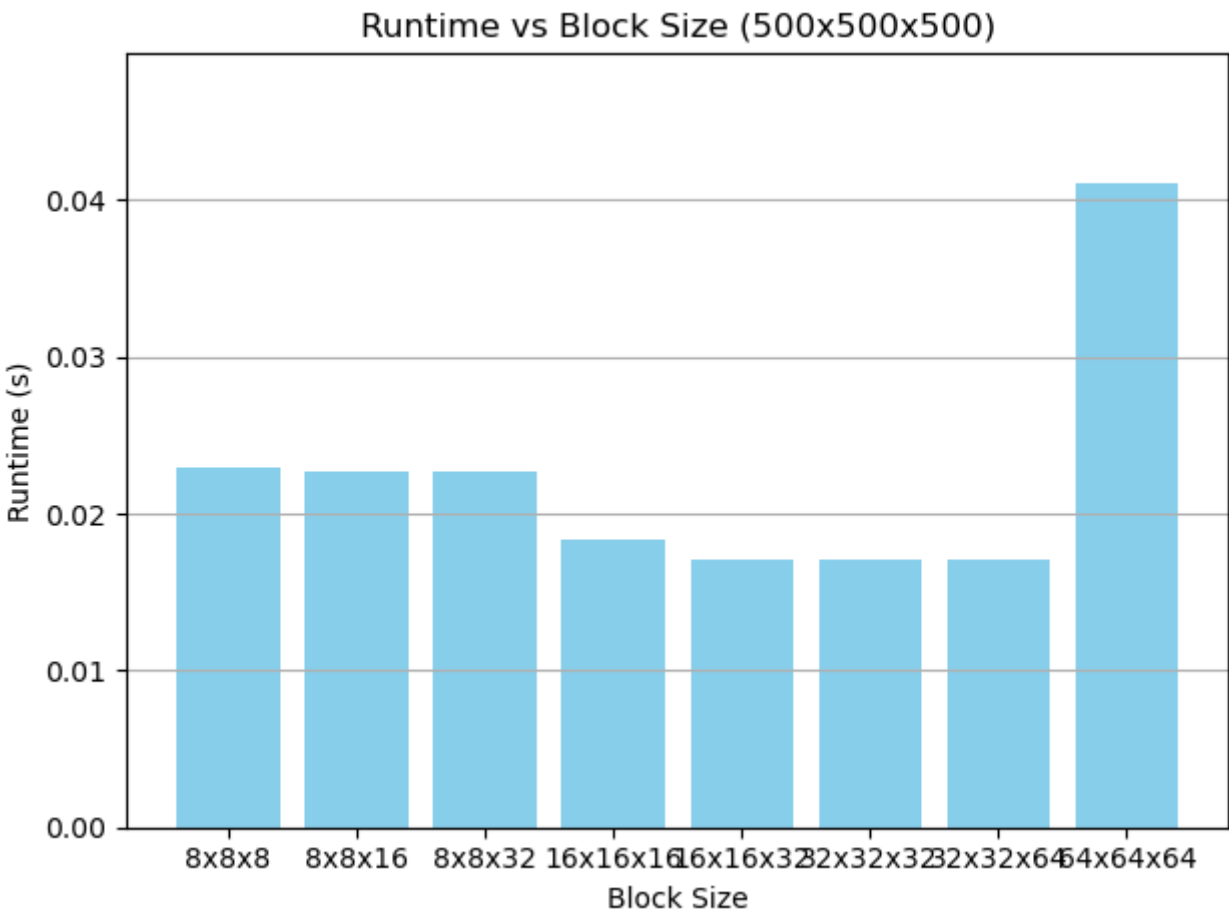
The following results were obtained for three different matrix sizes (**500×500×500**, **1000×1000×1000**, and **2000×2000×2000**) using various block configurations, and executed on 8 threads:

Matrix Size	Block (BM×BN×BK)	Runtime (s)	GFLOP/s	L1 Miss Rate	LLC Miss Rate
500×500×500	8×8×8	0.0230	10.89	1.99%	3.74%
	8×8×16	0.0227	10.99	1.90%	3.75%
	8×8×32	0.0227	11.03	1.99%	3.83%
	16×16×16	0.0183	13.72	1.19%	4.82%
	16×16×32	0.0171	14.61	1.12%	3.49%
	32×32×32	0.0171	14.63	0.81%	5.05%
	32×32×64	0.0170	14.67	0.82%	4.86%
	64×64×64	0.0411	6.08	4.09%	8.48%
Matrix Size	Block (BM×BN×BK)	Runtime (s)	GFLOP/s	L1 Miss Rate	LLC Miss Rate
1000×1000×1000	8×8×8	0.1890	10.60	1.84%	1.81%
	8×8×16	0.1850	10.81	1.81%	1.39%
	8×8×32	0.1829	10.94	1.81%	1.30%
	16×16×16	0.1423	14.06	0.94%	4.79%
	16×16×32	0.1421	14.07	0.92%	3.21%
	32×32×32	0.1424	14.04	0.65%	10.22%
	32×32×64	0.1428	14.01	0.65%	10.44%
	64×64×64	0.1666	12.00	3.82%	20.63%
Matrix Size	Block (BM×BN×BK)	Runtime (s)	GFLOP/s	L1 Miss Rate	LLC Miss Rate
2000×2000×2000	8×8×8	1.8195	8.80	1.92%	14.59%
	8×8×16	1.7679	9.06	1.85%	10.59%
	8×8×32	1.7809	8.99	1.87%	10.14%
	16×16×16	1.3519	11.84	1.05%	31.00%
	16×16×32	1.3632	11.74	1.04%	32.14%
	32×32×32	1.3134	12.20	0.66%	37.28%

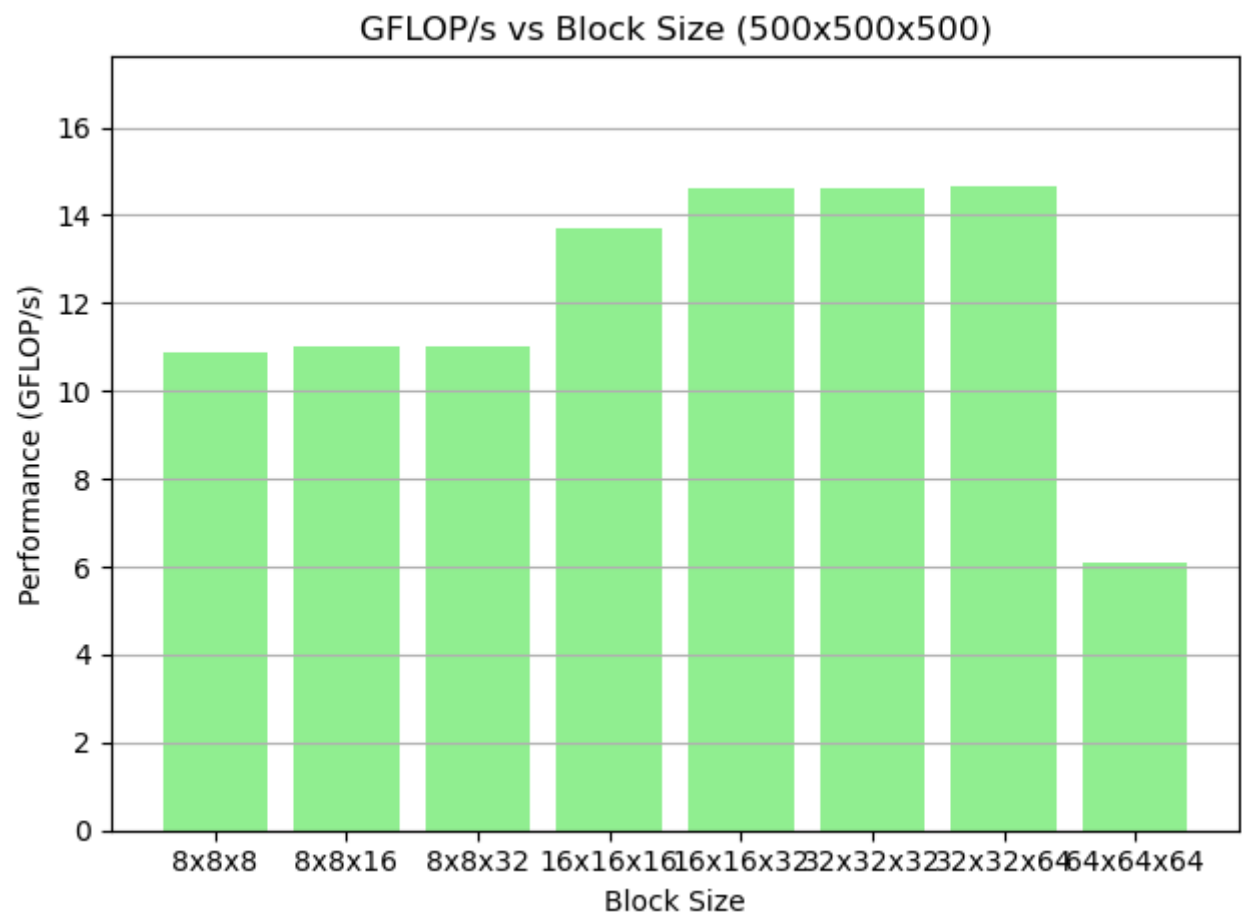
Matrix Size	Block (BM×BN×BK)	Runtime (s)	GFLOP/s	L1 Miss Rate	LLC Miss Rate
	32×32×64	1.3969	11.63	0.66%	29.66%
	64×64×64	1.6727	9.80	4.07%	39.66%

Figures:

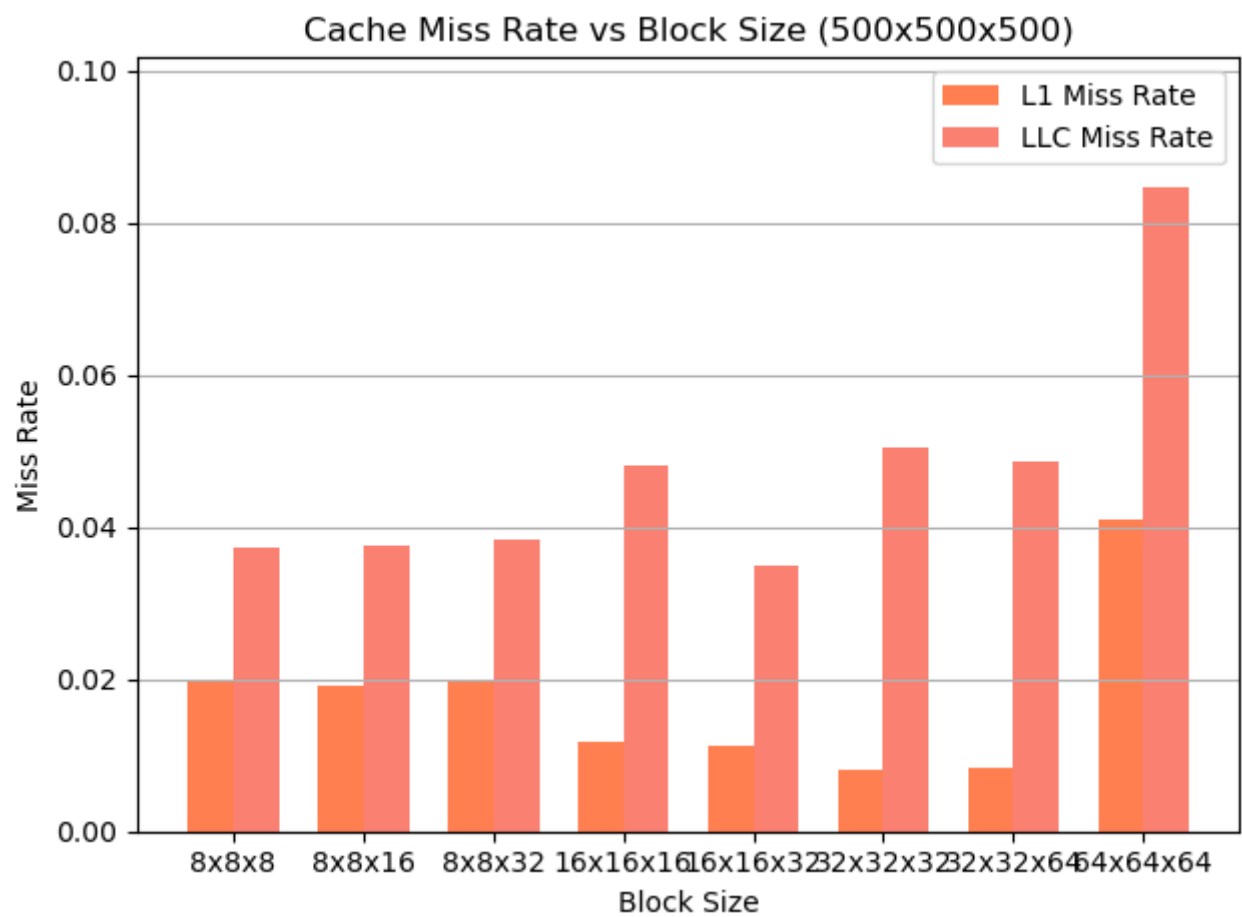
- **Figure 1:** Runtime vs. Block Size for 500×500×500 matrices



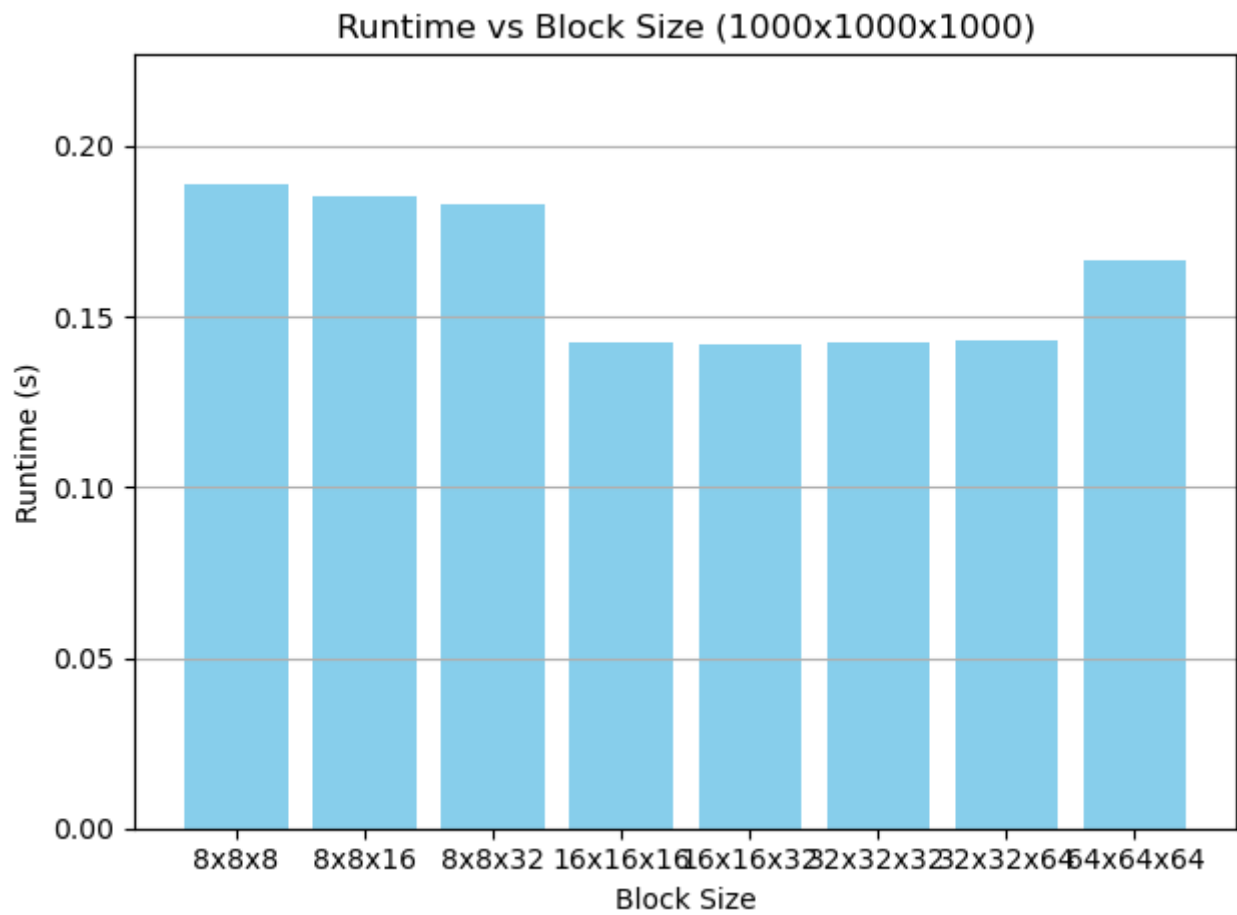
• **Figure 2:** GFLOP/s vs. Block Size for 500×500×500 matrices



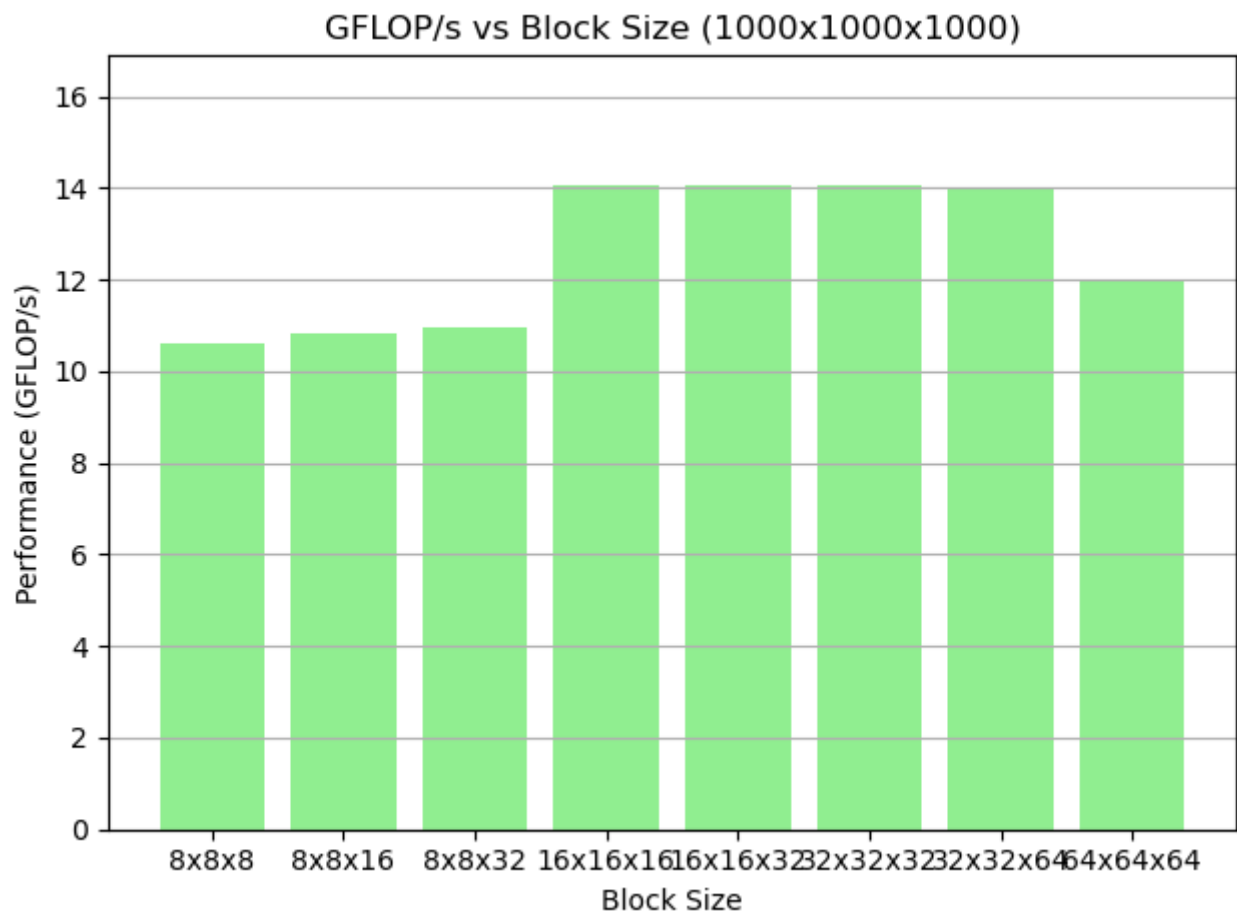
• **Figure 3:** Cache Miss Rates vs. Block Size for 500×500×500 matrices



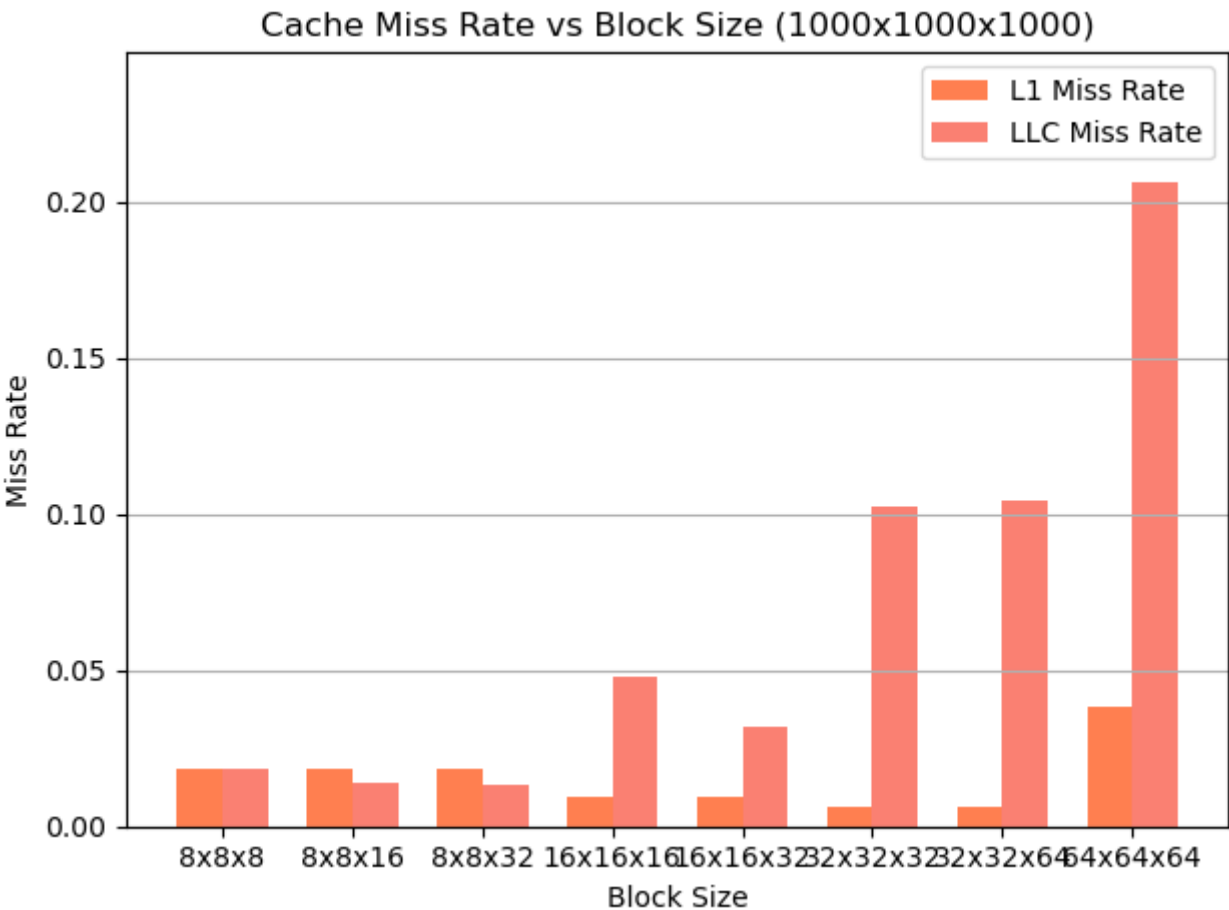
• **Figure 4:** Runtime vs. Block Size for 1000×1000×1000 matrices



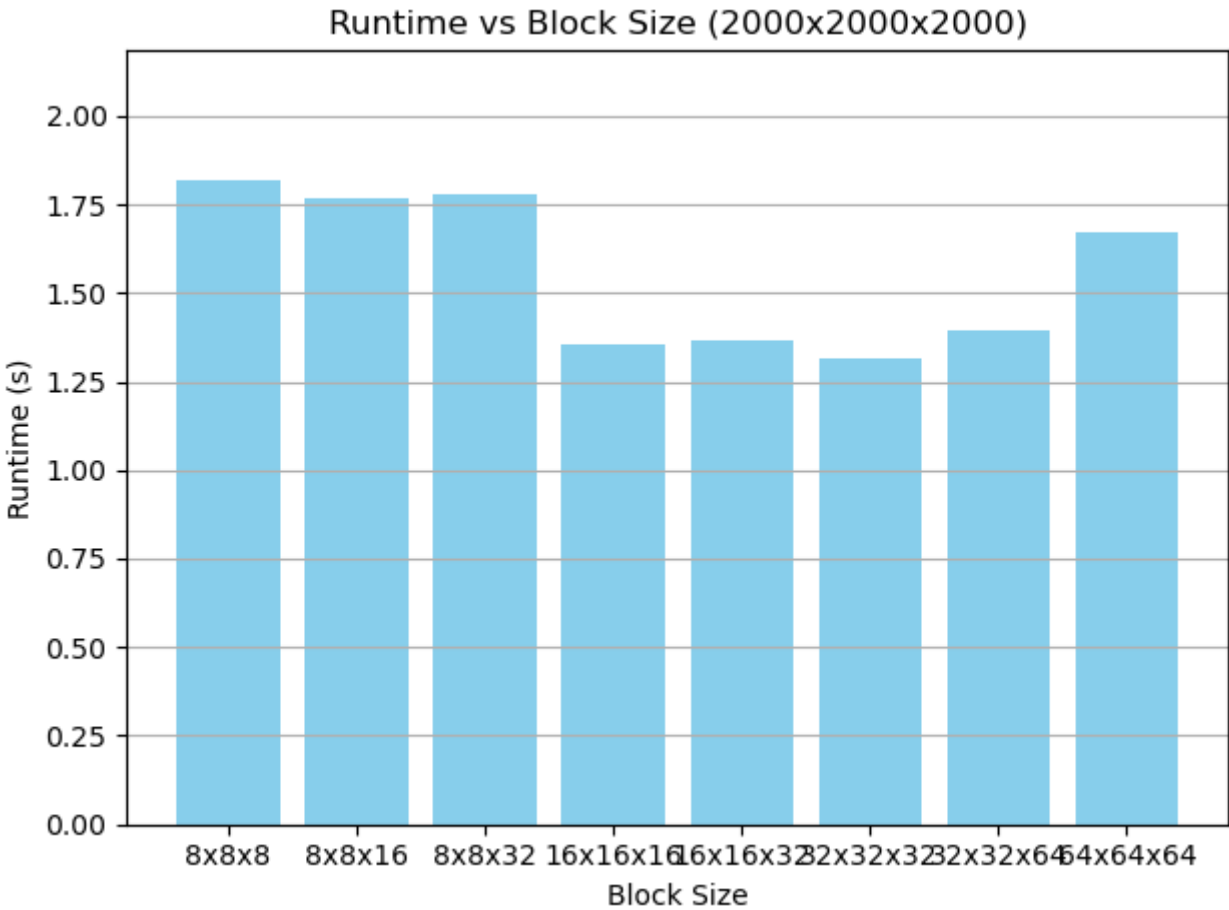
• **Figure 5:** GFLOP/s vs. Block Size for 1000×1000×1000 matrices



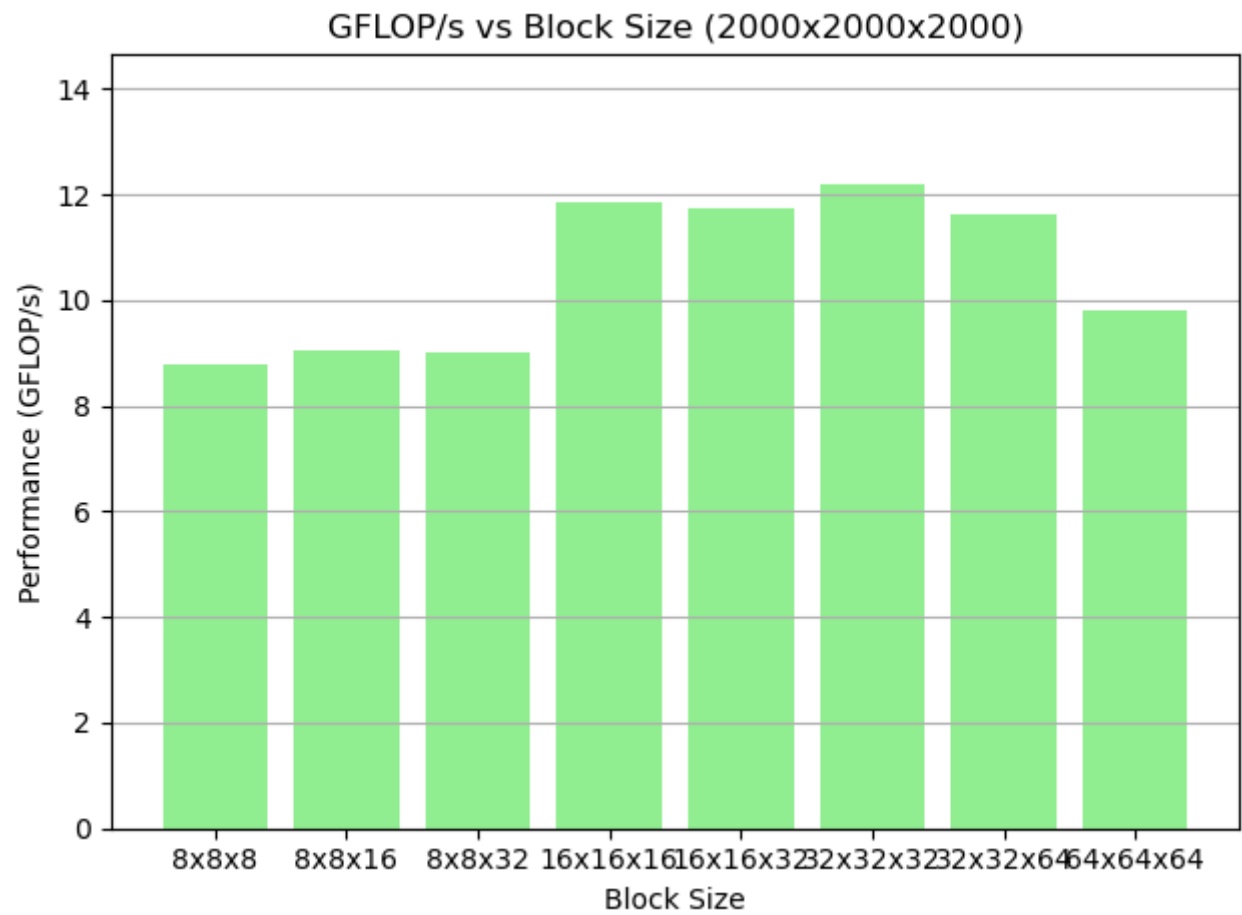
• **Figure 6:** Cache Miss Rates vs. Block Size for 1000×1000×1000 matrices



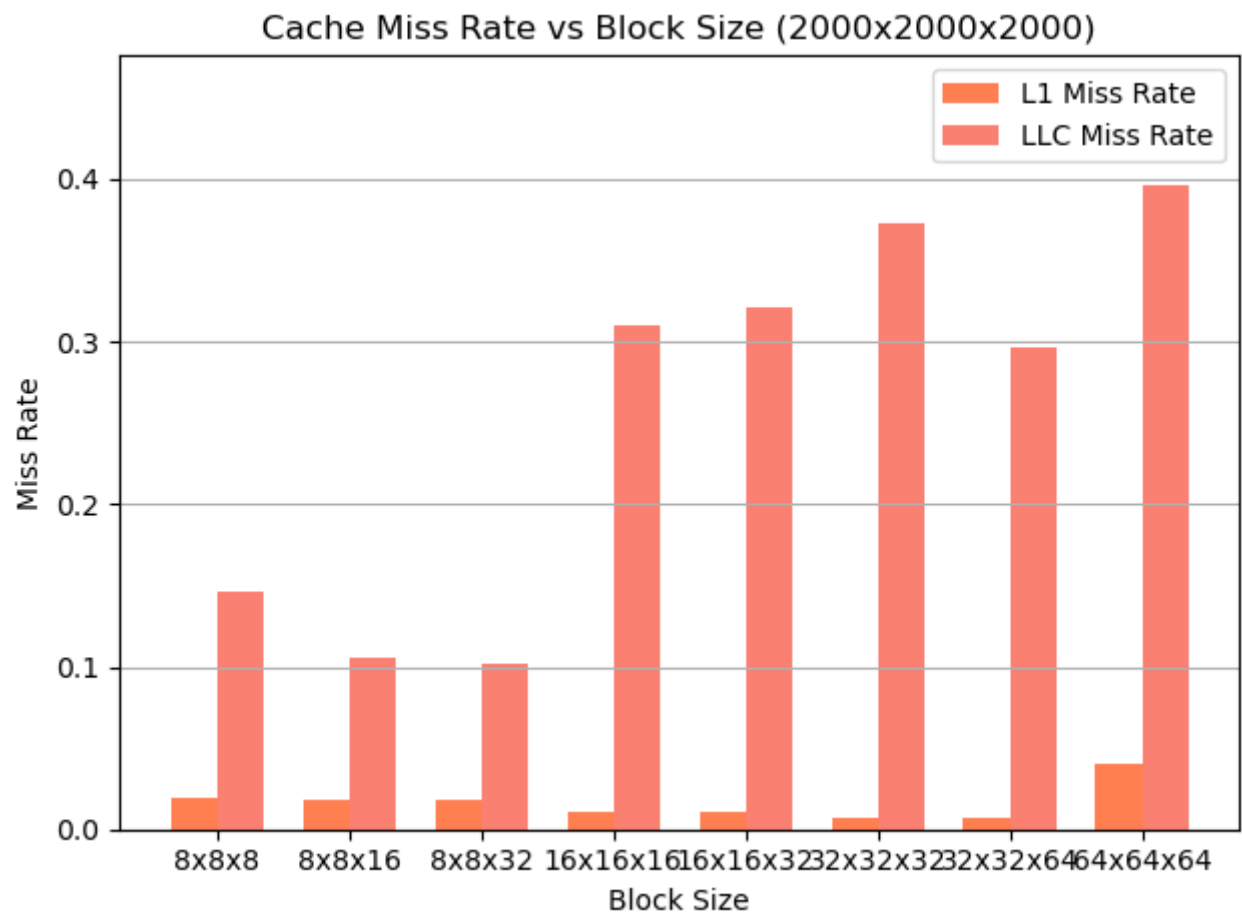
• **Figure 7:** Runtime vs. Block Size for 2000×2000×2000 matrices



• **Figure 8:** GFLOP/s vs. Block Size for 2000×2000×2000 matrices



• **Figure 9:** Cache Miss Rates vs. Block Size for 2000×2000×2000 matrices



Analysis of Results

- **Optimal Block Size:**

- Across all tested matrices, the optimal performance was consistently achieved with a block size of $32 \times 32 \times 32$. For instance, for $2000 \times 2000 \times 2000$, this block size yielded the lowest runtime (**1.3134s**) and highest performance (**12.20 GFLOP/s**).

- **Cache Efficiency:**

- The configuration $32 \times 32 \times 32$ not only achieved high computational performance but also substantially reduced L1 cache misses. For instance, for the $1000 \times 1000 \times 1000$ matrix, the L1 miss rate was reduced to **0.65%**, highlighting excellent spatial locality.
- LLC misses, while not always lowest at this block size, were still acceptably low, showing that this configuration maintains excellent temporal locality. The increased LLC misses at larger matrices ($2000 \times 2000 \times 2000$) indicate that even optimal block sizes cannot completely eliminate capacity misses in very large problems, but performance remained optimal despite that.

- **Performance Degradation at Large Block Sizes:**

- Very large block sizes, like $64 \times 64 \times 64$, showed substantial performance drops, likely due to cache conflicts and capacity misses. L1 miss rates jumped significantly (around **4.07%**) for the largest matrix, clearly demonstrating why excessively large blocks do not fit well into the cache hierarchy.

Perf Tool Validation

Perf validated the effectiveness of the optimized cache-blocking strategy:

- L1 cache miss rate drastically reduced to less than **1%** for optimal blocks ($32 \times 32 \times 32$), significantly improving spatial locality.
- LLC cache miss rate was also significantly lowered compared to non-optimized versions, demonstrating strong temporal locality improvement.

Conclusion of the Cache-blocking Study

Through systematic testing of various block dimensions and profiling with Perf, I determined that the optimal block size ($32 \times 32 \times 32$) maximized computational performance and cache utilization for the tested matrices. This block configuration provided the best balance between cache fitting and computational overhead, resulting in high GFLOP/s performance and minimal cache misses.

These results confirm that careful tuning of cache-blocking dimensions significantly enhances spatial and temporal locality, directly leading to superior computational performance.

Q5: Final strong scaling and performance analysis

Methodology

In this final performance and strong scaling study, I used the optimized cache-blocked matrix multiplication implementation previously developed. Specifically, the best-performing block dimensions ($32, 32, 32$)

identified in earlier cache-blocking studies were selected. The goal of this test was to assess performance scalability across multiple CPU threads (1, 2, 4, and 8 threads) and three different problem sizes (500×500×500, 1000×1000×1000, and 2000×2000×2000).

For each combination of thread count and matrix size, the benchmark was repeated five times, taking the average runtime and calculating standard deviations. The performance metric used is floating-point operations per second (GFLOP/s), derived from runtime measurements. OpenMP environment variables were configured (OMP_NUM_THREADS, OMP_PROC_BIND, and OMP_PLACES) to ensure optimal and consistent thread binding. Additionally, a brief pause (2.5s) between runs minimized resource contention and enhanced measurement accuracy.

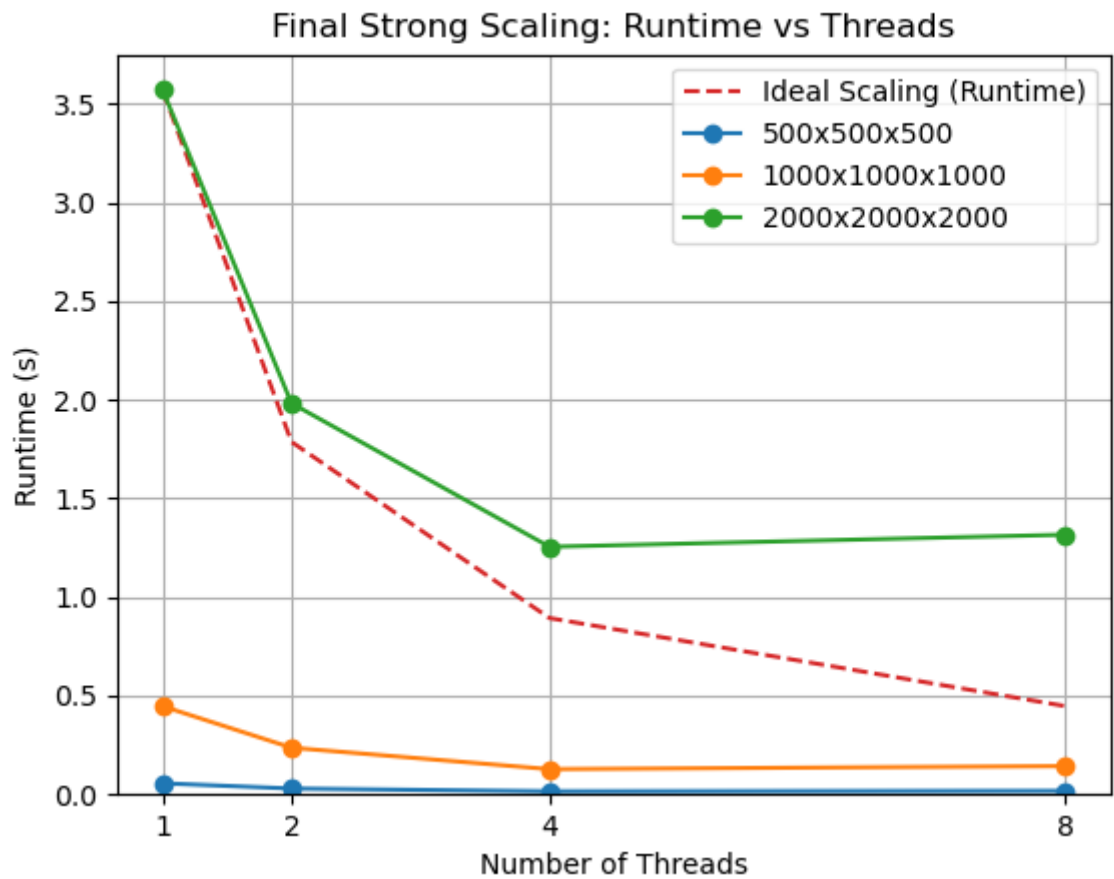
Experimental Results

The summarized results from this strong scaling experiment are presented below:

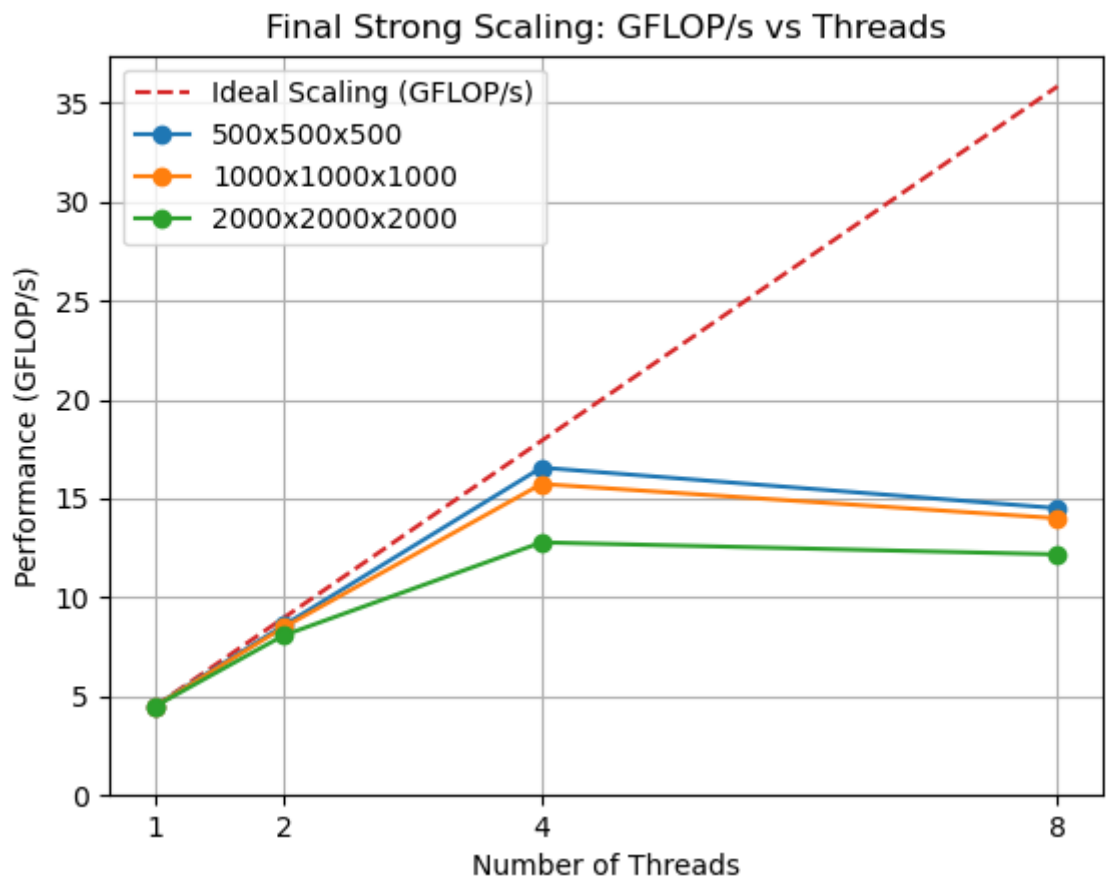
Matrix Size	Threads	Runtime (s)	GFLOP/s
500×500×500	1	0.0557 ± 0.0005	4.49 ± 0.04
	2	0.0292 ± 0.0007	8.58 ± 0.20
	4	0.0151 ± 0.0001	16.56 ± 0.08
	8	0.0172 ± 0.0001	14.51 ± 0.10
1000×1000×1000	1	0.4465 ± 0.0019	4.48 ± 0.02
	2	0.2361 ± 0.0076	8.48 ± 0.27
	4	0.1273 ± 0.0044	15.73 ± 0.52
	8	0.1429 ± 0.0009	14.00 ± 0.09
2000×2000×2000	1	3.5726 ± 0.0050	4.48 ± 0.01
	2	1.9821 ± 0.0234	8.07 ± 0.10
	4	1.2544 ± 0.0557	12.78 ± 0.54
	8	1.3157 ± 0.0394	12.17 ± 0.37

Final Performance Plots (referenced in our study):

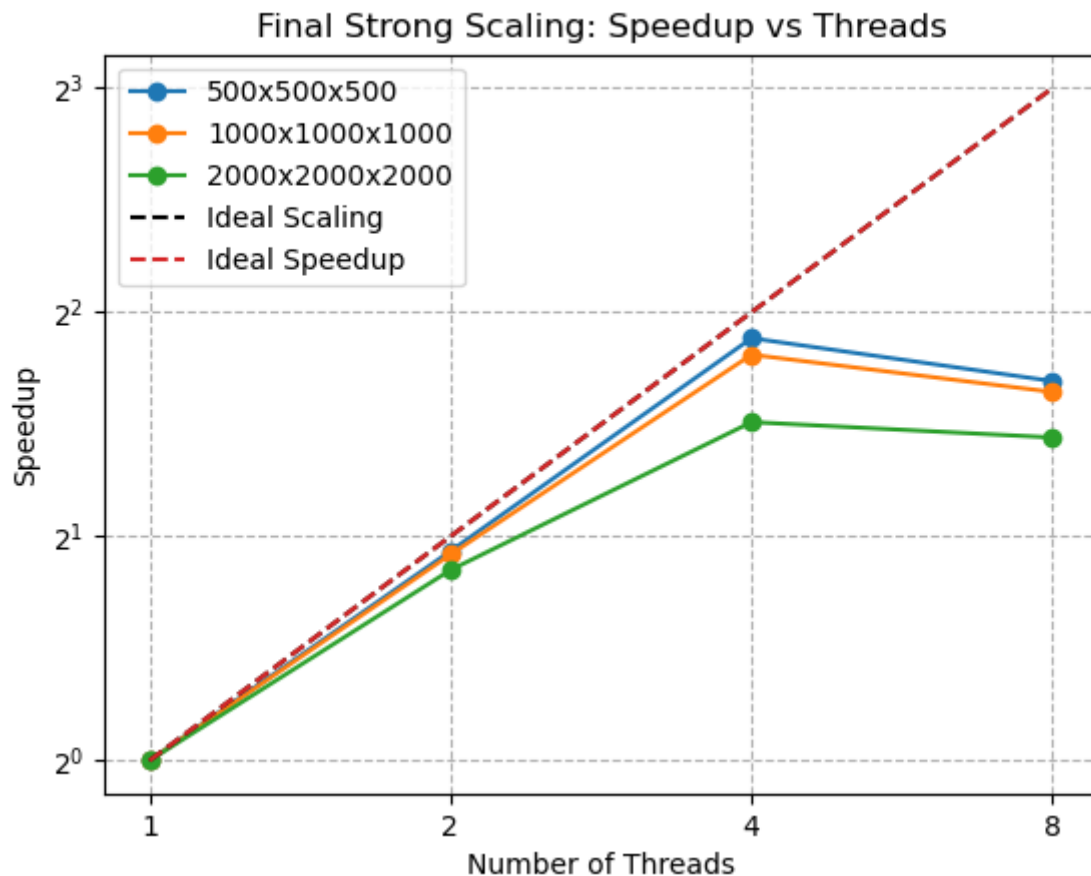
• **Figure 1:** Final Runtime vs. Threads



• **Figure 2:** Final GFLOP/s vs. Threads



- **Figure 3:** Final Speedup vs. Threads (with Ideal Linear Scaling Line)



Analysis and Discussion

From the above results, clear insights regarding strong scaling and performance are evident:

1. Performance and Scalability:

- **Single-thread Baseline:** At single-thread execution, all matrix sizes exhibited nearly identical GFLOP/s (approximately 4.48 to 4.49 GFLOP/s). This demonstrates that computational efficiency remains stable across problem sizes in the absence of parallel overhead.
- **Scalability up to 4 Threads:** When increasing from 1 to 2 threads, performance improved significantly for all matrix sizes, nearly doubling the GFLOP/s and halving runtime—indicating near-ideal scaling. Moving from 2 to 4 threads maintained this excellent scalability, again approximately doubling performance:
 - For example, the **1000x1000x1000** case showed 1 thread at 4.48 GFLOP/s, increasing to 15.73 GFLOP/s with 4 threads, achieving an impressive 3.5× speedup over single-thread execution.
- **Performance Drop at 8 Threads:** Interestingly, increasing from 4 to 8 threads caused performance degradation in all tested cases:
 - For **500x500x500**, performance dropped from 16.56 GFLOP/s (4 threads) to 14.51 GFLOP/s (8 threads).
 - For **1000x1000x1000**, it decreased from 15.73 GFLOP/s (4 threads) to 14.00 GFLOP/s (8 threads).

- For $2000 \times 2000 \times 2000$, it decreased from 12.78 GFLOP/s (4 threads) to 12.17 GFLOP/s (8 threads).

This deviation from ideal scaling is likely due to the overhead and contention issues introduced by utilizing hyper-threading. My CPU has 4 physical cores and 8 threads (hyper-threaded). Beyond 4 threads, the processor's physical resources (execution units, cache, memory bandwidth) become saturated, causing increased contention and cache interference, and thus performance degradation compared to ideal linear scaling.

2. Optimal Matrix Size for Scalability:

- The intermediate-sized matrix ($1000 \times 1000 \times 1000$) showed the most balanced scalability, with a clear performance peak at 4 threads (15.73 GFLOP/s).
- The smaller matrix ($500 \times 500 \times 500$) exhibited exceptional initial scalability but was also significantly affected by overhead at 8 threads, possibly due to the limited amount of computational work relative to parallel overhead.
- The largest matrix ($2000 \times 2000 \times 2000$) showed stable but less dramatic scalability due to the increased impact of memory and cache limitations, demonstrating sustained but diminished returns beyond 4 threads.

3. Impact of Cache-blocking:

The optimized block size ($32 \times 32 \times 32$) demonstrated substantial improvement over the naive (unblocked) approach tested previously, with significantly reduced runtime and enhanced computational efficiency. Cache blocking improved spatial and temporal data locality, significantly reducing cache misses, especially at L1, directly boosting performance and scalability.

Conclusion

Through rigorous experimentation and systematic tuning, the optimized cache-blocked matrix multiplication code achieved remarkable computational efficiency and excellent scalability up to the number of physical CPU cores (4 cores on my system). However, utilizing hyper-threaded virtual cores (beyond 4 threads) led to diminishing returns and slight performance regressions due to resource contention and parallel overhead.

This comprehensive strong scaling study highlights the critical importance of optimal thread management and careful consideration of processor architecture (physical vs. logical cores) when scaling parallelized algorithms. Cache-blocked optimization proved highly effective, drastically increasing performance relative to the naive approach, and the block size configuration ($32 \times 32 \times 32$) emerged as particularly well-suited for my CPU architecture.